

# LOCAL LINEAR REGRESSION FOR GENERALIZED LINEAR MODELS WITH MISSING DATA

C. Y. WANG<sup>1</sup>, Suojin WANG<sup>2</sup> R. J. CARROLL<sup>3</sup> and Roberto G. GUTIERREZ

*Fred Hutchinson Cancer Research Center, Texas A&M University, Texas A&M University and  
Southern Methodist University*

January 14, 1997

## Abstract

Fan, Heckman and Wand (1995) proposed locally weighted kernel polynomial regression methods for generalized linear models and quasilielihood functions. When the covariate variables are missing at random, we propose a weighted estimator based on the inverse selection probability weights. Distribution theory is derived when the selection probabilities are estimated nonparametrically. We show that the asymptotic variance of the resulting nonparametric estimator of the mean function in the main regression model is the same as that when the selection probabilities are known, while the biases are generally different. This is different from results in parametric problems, where it is known that estimating weights actually decreases asymptotic variance. To reconcile the difference between the parametric and nonparametric problems, we obtain a second-order variance result for the nonparametric case. We generalize this result to local estimating equations. Finite sample performance is examined via simulation studies. The proposed method is demonstrated via an analysis of data from a case-control study.

*Short title.* Local regression with missing data

---

<sup>1</sup>Research supported by National Cancer Institute grant (CA-53996).

<sup>2</sup>Research supported by grants from the National Science Foundation (DMS-9504589), the National Security Agency (MDA904-96-1-0029), the National Cancer Institute (CA-57030), and Texas A&M University's Scholarly and Creative Activities Program (95-59).

<sup>3</sup>Research supported by National Cancer Institute grant (CA-57030) and partially completed while visiting the Institut für Statistik und Ökonometrie, Sonderforschungsbereich 373, Humboldt Universität zu Berlin, with partial support from a senior Alexander von Humboldt Foundation research award.

*AMS 1991 subject classifications.* Primary 62G07; secondary 62G20

*Key words and phrases.* Generalized linear models, kernel regression, local linear smoother, measurement error, missing at random, quasilielihood functions.

# 1 INTRODUCTION

This paper is concerned with nonparametric function estimation via quasilikelihood when the predictor variable may be missing, and the missingness depends upon the response. We use local polynomials with kernel weights, generalizing the work of Staniswalis (1989), Severini and Staniswalis (1994) and Fan, Heckman and Wand (1995) to the missing data problem.

In practice, covariates may be missing due to reasons such as loss to follow up. For example, in a study of acute graft versus host disease of bone marrow transplants of 97 female subjects conducted at the Fred Hutchinson Cancer Research Center, the outcome is the acute graft host disease and one covariate of interest is the donor's previous pregnancy status which was missing for 31 patients because of the incompleteness of the donors' medical history. In this paper, we consider the missing covariate data problem in nonparametric generalized linear models. We assume that covariates are missing at random (MAR) and the missingness is ignorable (Rubin, 1976).

In parametric problems, two approaches are common. Likelihood methods assume a joint parametric distribution for covariates and response, and under our assumptions ignore the missing data mechanism (Little and Rubin, 1987). Complete-case analysis assumes nothing about the distribution of covariates, and is in this sense semiparametric. Estimation is based on the "complete-cases", *i.e.*, those with no missing data, with weighting inversely proportional to the probability that the covariate is observed given the response (Horvitz and Thompson, 1952). We call these *selection probabilities*. We use the second approach. Our methods apply as well to other semiparametric schemes, *e.g.*, that of Robins, Rotnitzky and Zhao (1994). We estimate the missing data probabilities by nonparametric regression.

In parametric problems, the Horvitz-Thompson weighting scheme has a curious and important property. Consider two estimators: (a) the one with known selection probabilities and weights; and (b) one where the selection probabilities are estimated by a properly specified parametric model. The two methods yield consistent estimates, but that with estimated weights generally has a *smaller* asymptotic variance (Robins, et al., 1994).

One might expect the same sort of result to hold in the nonparametric regression case with nonparametrically estimated selection probabilities. However, this is not the case, and we show (Theorem 1) that whether weights are estimated or not has no effect on asymptotic variance, while it does have an effect on the bias in general.

In simulations however, we observed repeatedly that estimating weights was beneficial in small samples. To understand whether this numerical evidence was at all general, we developed a second-order variance result (Theorem 2) showing that the estimator with estimated weights can be expected to have smaller finite-sample variance than if the weights are known. This second-order variance result provides a reconciliation between the different first-order results in the parametric and nonparametric cases.

The statistical models are described in Section 2. In Section 3, we propose the methodology and the asymptotic result for the weighted method with both known and estimated selection probabilities. The method is demonstrated in Section 4 by analyzing the data from a case-control study of bladder cancer. In Section 5 we investigate the finite sample performance by conducting a simulation study. We note that estimating the selection probabilities has a finite sample effect on the estimation of the mean function of our primary interest. We explain the possible finite sample efficiency gain by a second-order variance approximation in Section 6.

The major result of Section 3 can be described as follows:

- An unknown function  $\pi(\cdot)$  is estimated nonparametrically, by  $\hat{\pi}(\cdot)$ .
- If  $\pi(\cdot)$  were known, one would use it to estimate nonparametrically a second function  $\mu(\cdot)$ , by  $\hat{\mu}(\cdot, \pi)$ .
- The estimates  $\hat{\mu}(\cdot, \pi)$  and  $\hat{\mu}(\cdot, \hat{\pi})$  have the same asymptotic variance.

In Section 7, we sketch a result showing that this phenomenon is quite general, and not restricted to our particular context. All detailed proofs are given in the Appendix.

## 2 THE MODELS

### 2.1 Full Data Models

We let  $(Y_1, X_1), \dots, (Y_n, X_n)$  be a set of independent random variables, where  $Y_i$  is a scalar response variable, and  $X_i$  is a scalar covariate variable. In a classical generalized linear model (Nelder and Wedderburn, 1972; McCullagh and Nelder, 1989), the conditional density of  $Y$  given  $X$  belongs to a canonical exponential family  $f_{Y|X}(y|x) = \mathcal{C}(y)\exp[y\theta(x) - \mathcal{B}\{\theta(x)\}]$  for known functions  $\mathcal{B}$  and  $\mathcal{C}$ , where the function  $\theta$  is called the canonical or natural parameter. The unknown function  $\mu(x) = E(Y|X = x)$  is modeled in  $X$  by a link function  $g$  by  $g\{\mu(x)\} = \eta(x)$ . In a parametric

generalized linear model,  $\eta(x) = c_0 + c_1x$  for some unknown parameter  $c_0, c_1$ . The link function  $g$  is assumed to be known. For example, in logistic regression  $g(u) = \log\{u/(1-u)\}$ , and in linear regression  $g(u) = u$ . In our nonparametric setting, there is no model assumption about  $\eta(x)$ .

Fan, et al. (1995) considered quasilielihood models, where only the relationship between the mean and the variance is specified. If the conditional variance is modeled as  $\text{var}(Y|X = x) = V\{\mu(x)\}$ , for some known positive function  $V$  then the corresponding quasilielihood function  $Q(w, y)$  satisfies  $(\partial/\partial w)Q(w, y) = (y - w)/V(w)$  (Wedderburn, 1974). The primary interest is to estimate  $\mu(x)$ , or equivalently  $\eta(x)$ , nonparametrically.

## 2.2 Missing Data Models

In a missing covariate data problem, some covariates may be missing and we let  $\delta_i = 1$  if  $X_i$  is observed,  $\delta_i = 0$  otherwise. Furthermore, let

$$\pi_i = \text{pr}(\delta_i = 1|Y_i, X_i) = \text{pr}(\delta_i = 1|Y_i) = \pi(Y_i) \quad (1)$$

be the selection probability which does not depend on  $X_i$ , i.e.,  $X_i$  is MAR. In a two-stage design (White, 1982), often the selection probabilities are known. In many missing data problems, however, the selection probabilities are unknown and need to be estimated. To model the selection probabilities, we assume that given  $Y$  there is a known link function  $g^*$  such that  $g^*\{\pi(y)\} = \eta^*(y)$ , where  $\eta^*(y)$  is a smooth function. Let the conditional variance be modeled by  $\text{var}(\delta|Y = y) = V^*\{\pi(y)\}$  for some known positive function  $V^*$ . The corresponding quasilielihood function  $Q^*(w, \delta)$  satisfies  $(\partial/\partial w)Q^*(w, \delta) = (\delta - w)/V^*(w)$ . We say that *two-stage data models* occur when the selection probabilities are known, and *missing data models* occur when the selection probabilities are unknown. In the missing data models,  $\pi(y)$ , or  $\eta^*(y)$ , is a nuisance component which needs to be estimated.

## 3 METHODOLOGY

### 3.1 The Weighted Method

When  $(Y_i, X_i)$  are fully observable, Fan, et al. (1995) proposed the local linear kernel estimator of  $\eta(x)$  as  $\hat{\eta}(x, h) = \hat{\beta}_0$ , where  $h$  is the bandwidth of a kernel function  $K$  and  $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1)^t$  maximizes

$$\sum_{i=1}^n Q[g^{-1}\{\beta_0 + \beta_1(X_i - x)\}, Y_i] K_h(X_i - x), \quad (2)$$

where  $K_h(\cdot) = K(\cdot/h)$ . We assume that the maximizer exists, and this can be verified for standard choices of  $Q$ . The mean function  $\mu(x)$  is estimated by  $\hat{\mu}(x) = g^{-1}(\hat{\beta}_0)$ . When data are missing, a naive method is to apply (2) by using the complete-case (CC) analysis, *i.e.*, solving (2) by restricting to pairs in which both  $Y$  and  $X$  are observed. However, complete-case analysis may cause considerable bias when the missingness probabilities (1) depend on the response (Little and Rubin, 1987).

To accommodate the missingness in the observed data, we propose a Horvitz-Thompson inverse-selection weighted method, so that the estimator of  $\beta$  maximizes

$$\sum_{i=1}^n Q[g^{-1}\{\beta_0 + \beta_1(X_i - x)\}, Y_i] \frac{\delta_i}{\pi(Y_i)} K_h(X_i - x). \quad (3)$$

Note that here  $\pi(Y_i)$  is assumed to be known and strictly positive in the support of  $Y$ . For notational purposes, we denote the solution to (3) by  $\hat{\beta}(\pi)$ .

We now define some notation for the presentation of the asymptotic properties of  $\hat{\beta}_0 = \hat{\eta}(x, \pi)$ . Suppose that  $K$  is supported on  $[-1, 1]$ . For any set  $\mathcal{A} \subset R$ , and  $i = 0, 1, 2, 3$ , let  $\gamma_i(\mathcal{A}) = \int_{\mathcal{A}} z^i K(z) dz$ ,  $\tau_i(\mathcal{A}) = \int_{\mathcal{A}} z^i K^2(z) dz$ . Define

$$\begin{aligned} N_x^h &= \{z : x - hz \in \text{supp}(f_X)\} \cap [-1, 1], \\ b_x &= \frac{1}{2} \eta^{(2)}(x) [g'\{\mu(x)\}]^{-1} \frac{\gamma_2^2(N_x^h) - \gamma_1(N_x^h) \gamma_3(N_x^h)}{\gamma_0(N_x^h) \gamma_2(N_x^h) - \gamma_1^2(N_x^h)}, \\ \sigma_x^2 &= f_X^{-1}(x) \mathcal{L}(x) \frac{\gamma_2^2(N_x^h) \tau_0(N_x^h) - 2\gamma_1(N_x^h) \gamma_2(N_x^h) \tau_1(N_x^h) + \gamma_1^2(N_x^h) \tau_2(N_x^h)}{\{\gamma_0(N_x^h) \gamma_2(N_x^h) - \gamma_1^2(N_x^h)\}^2}, \end{aligned}$$

where  $f_X(x)$  is the density of  $X$  and

$$\mathcal{L}(x) = E \left[ \frac{\{Y_1 - \mu(x)\}^2}{\pi(Y_1)} \middle| X_1 = x \right]. \quad (4)$$

As we will see later,  $\sigma_x^2$  is the asymptotic variance of  $\hat{\mu}(x, \pi)$ . For a bandwidth  $h$ ,  $x$  is an interior point of  $\text{supp}(f_X)$  if and only if  $N_x^h = N_{x,h} = [-1, 1]$ . To estimate  $\mu(x) = g^{-1}\{\eta(x)\}$ , we let  $\hat{\mu}(x, \pi) = g^{-1}\{\hat{\eta}(x, \pi)\} = g^{-1}(\hat{\beta}_0)$ . The limit distribution of  $\hat{\mu}(x, \pi)$  presented in Theorem 1 below can be obtained by calculations similar to that in Fan, et al. (1995).

### 3.2 Main Theorem

We now investigate the case with unknown selection probabilities. To estimate the selection probabilities, we again apply the local linear smoother of Fan, et al. (1995). For a fixed point  $y$ , we estimate  $\pi(y)$  by

$$\hat{\pi}(y) = g^{*-1}(\hat{\alpha}_0), \quad (5)$$

where  $\hat{\alpha} = (\hat{\alpha}_0, \hat{\alpha}_1)$  maximizes  $\sum_{i=1}^n Q^*[g^{*-1}\{\alpha_0 + \alpha_1(Y_i - y), \delta_i\}]K_\lambda(Y_i - y)$ , where we use  $\lambda$  as the smoothing parameter to distinguish it from the other smoothing parameter  $h$  used in estimating  $\beta$  for estimating the primary mean function  $\mu$ . Note that if the outcome  $Y$  is categorical such as the situation in Section 5.2 or the data analysis in Section 4, then as  $\lambda \rightarrow 0$  the estimate of  $\pi$  is equal to the empirical averages.

Let  $\hat{\beta}(\hat{\pi})$  maximize

$$\sum_{i=1}^n Q[g^{-1}\{\beta_0 + \beta_1(X_i - x)\}, Y_i] \frac{\delta_i}{\hat{\pi}(Y_i)} K_h(X_i - x), \quad (6)$$

where  $\hat{\pi}(y)$  is given in (5). Similar to the definition of  $\hat{\mu}(x, \pi)$ , we define  $\hat{\mu}(x, \hat{\pi}) = g^{-1}\{\hat{\eta}(x, \hat{\pi})\}$  where  $\hat{\eta}(x, \hat{\pi}) = \hat{\beta}_0(\hat{\pi})$ . We now present our main result.

**Theorem 1.** Suppose that Conditions (A1)-(A7), (B1)-(B6) in the Appendix, are satisfied. Then if  $h = h_n \rightarrow 0$ ,  $nh^3 \rightarrow \infty$ , and  $\lambda = \lambda_n = c^*h$  for a constant  $c^* > 0$ , we have that for any  $x \in \text{supp}(f_X)$ , there exist  $b_{nj}(x) = b_x\{1 + o(1)\}$ ,  $j = 1, 2$ , such that both  $(nh)^{1/2}\{\hat{\mu}(x, \pi) - \mu(x) - h^2b_{n1}(x)\}$  and  $(nh)^{1/2}\{\hat{\mu}(x, \hat{\pi}) - \mu(x) - h^2b_{n2}(x) - \lambda^2 f_X(x)S_3(x)\}$  converge in distribution to a normal random variable with mean 0 and variance  $\sigma_x^2$ , where  $S_3(x)$  is given in (22) in the Appendix, and  $S_3(x) = 0$  if either  $Y$  is a lattice random variable or  $\pi$  is a constant.

One important implication of this result is that the asymptotic effect on the asymptotic variance due to estimating selection probabilities, which is nonnegligible in the parametric or semiparametric models (Robins, et al., 1994; Wang, Wang, Zhao and Ou, 1997), disappears in the corresponding fully nonparametric problems. The difference appears in the bias term, but it vanishes if either  $Y$  is a lattice random variable or  $\pi$  is a constant. The proof of Theorem 1 is in the Appendix.

## 4 DATA ANALYSIS

In this section we consider an example of a case-control study of bladder cancer conducted at the Fred Hutchinson Cancer Research Center. Eligible subjects were residents of 3 counties of western Washington state who were diagnosed between January 1987 and June 1990 with invasive or noninvasive bladder cancer. This population based case-control study was designed to address the association between bladder cancer and some nutrients. We use the data here for illustrative purposes. Some detailed results can be found in Bruemmer, White, Vaughan and Cheney (1995).

In our demonstration, the response variable is the bladder cancer history and the covariate  $X$  is the smoking package year. The smoking package year of a participant is defined as the average

number of cigarette packages smoked per day multiplied by the years one has been smoking. There are a total of 262 cases and 405 controls. However, the smoking package year information of 1 case and 215 controls were missing. In addition, we treated past smokers as in the nonvalidation set since we are primarily interested in the smoking effect of current smokers. One case with  $X = 200$  has high leverage ( $X$  has mean 26 and standard deviation 30) and was not included in the validation set. As a result, there were 167 cases and 179 controls in the validation set.

To analyze the data, one may consider the complete-case logistic regression of  $Y$  on  $X$ , with and without adjustment by estimated inverse selection weights. The estimates of the slope (s.e.) are .0276 (.0047) and .0268 (.0046), respectively. The resulting estimates of  $E(Y|X)$ , called global estimates, are given in Figure 1. We note that a parametric estimator is based on global estimation. Based on this logistic regression analysis, one would argue that the risk of developing bladder cancer increases monotonically as a function of the average smoking year.

Alternatively, we may employ the weighted local estimation method. We used the Epanechnikov kernel function that  $K(u) = .75(1 - u^2)$  on  $[-1, 1]$ . The unweighted estimates of  $E(Y|X)$ , denoted by  $\hat{\mu}_{CC}(\cdot)$ , and the weighted estimator,  $\hat{\mu}(\cdot, \hat{\pi})$ , are given in Figure 1. Based on the bandwidth selection criteria given in the Appendix, we used 24.2 as the bandwidth for the weighted local smoother and 19.6 for the unweighted one. We notice that the CC analysis has basically captured the effect of the average package year, as it is somewhat parallel to  $\hat{\mu}(\cdot, \hat{\pi})$ . Based on this nonparametric analysis, the argument is somewhat different from the previous parametric one. For example, the curves between  $X = 40$  and  $X = 95$  do not increase as much as the other two segments ( $X < 40$  or  $X > 95$ ). Although it is true that the average package year has a significant effect on bladder cancer, our analysis suggests that piecewise logistic regression is more proper if parametric inference is to be made.

One small point concerns the interpretation of Figure 1. Prentice and Pyke (1979) showed that in a case-control study with an ordinary parametric logistic regression model, the logits of the observed case-control data differ from that of the population only in the intercept term. The same is true in our problem. This means that the basic monotonicities and flatness observed in Figure 1 are not affected by the case-control sampling, although the levels of estimated disease probability of course would differ.

## 5 SIMULATION STUDIES

We conducted simulations to better understand the finite sample performance of the weighted estimator and the finite sample effect due to estimating the selection probabilities. Recall that  $\hat{\mu}_{CC}$  is the unweighted method which applies the local linear smoother of Fan, et al. (1995) directly to the validation set only. We compare the biases and variances of  $\hat{\mu}_{CC}$ ,  $\hat{\mu}(\cdot, \pi)$  and  $\hat{\mu}(\cdot, \hat{\pi})$ .

### 5.1 Continuous Response

In this subsection we consider the case of continuous response  $Y$ . First we generated  $n = 200$   $X$ 's from a uniform  $[-1,1]$  distribution and the response variable  $Y$ 's follow the linear link such that  $Y_i = \mu(X_i) + .3\epsilon_i$ , where  $\mu(x_i) = x_i^2$ ,  $\epsilon_i$  ( $i = 1, \dots, n$ ) is a random sample from normal  $(0,1)$  distribution and are independent of  $X_i$ . The selection probability given  $Y$  is from the logistic model with intercept 0.0 and slope 1.0. Approximately 42% of the data are missing under the above selection probabilities. We ran 1,000 independent replicates in this simulation experiment, and we applied the linear link and logit link to estimate  $\mu(\cdot)$  and  $\pi(\cdot)$ , respectively. In each replicate,  $\hat{\mu}_{CC}(\cdot)$ ,  $\hat{\mu}(\cdot, \pi)$  and  $\hat{\mu}(\cdot, \hat{\pi})$  were obtained using the Epanechnikov kernel function  $K(u) = .75(1 - u^2)$  on  $[-1,1]$  and the data-driven bandwidth selection criteria as described in the Appendix.

The empirical biases of the estimators are shown in Figure 2 for  $x \in (-1,1)$ . The curves are the averages of the bias estimates over 1,000 runs. Note that the CC analysis has considerable bias and that  $\hat{\mu}(\cdot, \pi)$  and  $\hat{\mu}(\cdot, \hat{\pi})$  are very close in most points. Figure 3 shows the sample variances of  $\hat{\mu}(x, \pi)$  and  $\hat{\mu}(x, \hat{\pi})$ . It appears that the weighted estimator using estimated selection probabilities is at least as efficient as the one using the true  $\pi(\cdot)$ . There is considerable gain using estimated  $\pi$  for a range of  $X$  values, especially when  $X$  is around zero. The relative efficiency of  $\hat{\mu}(x, \hat{\pi})$  to  $\hat{\mu}(x, \pi)$  at  $x = 0$  is 1.29 when  $n = 200$ . If we increase the sample size to  $n = 2,000$ , then the corresponding relative efficiency is 1.22. In Section 6, we explain the finite sample efficiency gain from estimating the selection probabilities by a second-order variance approximation.

### 5.2 Binary Response

We now study an important case when the response is binary. We generated  $n = 200$   $X$ 's from uniform  $[-1,1]$  distribution and the binary response  $Y$  was generated by

$$\text{pr}(Y_i = 1|X_i = x) = \mu(x) = \{1 + \exp(1 - x - x^2)\}^{-1}.$$



The selection probabilities depend on  $Y$  and are from a logistic model with intercept 1.0 and slope 1.0, leading to approximately 33% of the  $X$ 's being missing.

We now consider the nonparametric estimates. We applied the logit link for the estimation of both  $\mu(\cdot)$  and  $\pi(\cdot)$ . Because  $Y$  is binary,  $\pi(Y)$  was estimated by the empirical average at the corresponding  $Y$  value. The empirical biases of the resulting estimates  $\hat{\mu}(\cdot, \pi)$  and  $\hat{\mu}(\cdot, \hat{\pi})$  from 1,000 runs are again almost identical but the empirical variance of the latter is smaller. For  $n = 200$  the relative efficiency of  $\hat{\mu}(x, \hat{\pi})$  to  $\hat{\mu}(x, \pi)$  at  $x = .50$  is 1.21. When the sample size is increased to  $n = 2,000$ , the corresponding relative efficiency is 1.18. These findings are similar to the previous case with continuous response.

## 6 SECOND-ORDER VARIANCE APPROXIMATION

The simulations in the previous section show that there is finite sample gain from estimating the selection probabilities. Recall that the first-order asymptotic result of Theorem 1 shows no asymptotic efficiency gain from estimating the selection probabilities. To explain this, we now present the second-order variance approximation. The proof is given in the Appendix.

**Theorem 2.** Under the same conditions as in Theorem 1 and for any  $x \in \text{supp}(f_X)$  with  $\text{var}\{\hat{\mu}(x, \pi)\} < \infty$ , there exists  $\hat{\mu}_*(x) = \hat{\mu}(x, \hat{\pi}) + o_p(h^{1/2}n^{-1/2})$ , such that

$$\text{var}\{\hat{\mu}_*(x)\} = \text{var}\{\hat{\mu}(x, \pi)\} - n^{-1}v(x)\{1 + o(1)\},$$

for some  $v(x) > 0$ .

Theorem 2 shows that using the estimated selection probabilities gains the efficiency at the rate of  $n^{-1}$ . Note that the second-order efficiency gain is valid even when  $Y$  is a lattice random variable. For a fixed point  $x$ , let the relative efficiency gain by using the estimated selection probabilities be defined by  $[\text{var}\{\hat{\mu}(x, \pi)\} - \text{var}\{\hat{\mu}_*(x)\}]/\text{var}\{\hat{\mu}_*(x)\}$ . It is easy to see from Theorem 2 that the relative efficiency gain is of order  $O(h)$ , which goes to zero slowly. This supports the results of our simulations.

## 7 GENERALIZATIONS

Theorem 1 is a special case of a general phenomenon, which we outline here. Suppose that one has interest in a function  $\eta(\cdot)$ . If a nuisance function  $\pi(\cdot)$  were known, one would estimate  $\eta(\cdot)$  at  $x$  by

solving a local estimating equation of the form

$$0 = n^{-1} \sum_{i=1}^n K_h(X_i - x) \Psi\{\tilde{Y}_i, \pi(Z_i), \beta_0 + \beta_1(X_i - x)\} \{1, (X_i - x)\}^t, \quad (7)$$

where  $\Psi$  is an estimating function,  $Z$  is the covariate variable for  $\pi(\cdot)$  and  $\tilde{Y}$  represents a vector which may or may not include  $Z$ . In our problem, both  $\tilde{Y}$  and  $Z$  equal the response  $Y$ .

Now suppose that  $\pi(z)$  is also estimated by a local estimating equation but with bandwidth  $\lambda$ , so that

$$0 = n^{-1} \sum_{i=1}^n K_\lambda(Z_i - z) \Phi\{\tilde{Y}_i, X_i, \alpha_0 + \alpha_1(Z_i - z)\} \{1, (Z_i - z)\}^t.$$

The estimating functions  $\Psi$  and  $\Phi$  are assumed to satisfy

$$0 = E \left[ \Psi\{\tilde{Y}, \pi(Z), \eta(X)\} \right]; \quad 0 = E \left[ \Phi\{\tilde{Y}, X, \pi(Z)\} | Z \right].$$

Under this setup, in Appendix A.3 we sketch a result showing that

- The bias of  $\hat{\eta}(x)$  is of order  $h^2$ , is independent of the design densities of  $(Z, X)$ , but is generally affected by the estimation of  $\pi(\cdot)$ .
- The variance of  $\hat{\eta}(x)$  is asymptotically the same as if  $\pi(\cdot)$  were known.

Both these conclusions are reflected in our Theorem 1.

## ACKNOWLEDGEMENT

We are grateful to Barbara Bruemmer and Emily White for the case-control data.

## REFERENCES

- Bruemmer, B., White, E., Vaughan, T. and Cheney, C. (1995), “Nutrient Intake in Relationship to Bladder Cancer Among Middle aged Men and Women,” *Journal of National Cancer Institute*, in press.
- Carroll, R. J., Ruppert, D. and Welsh, A. H. (1996), “Nonparametric Estimation via Local Estimating Equations, with Applications to Nutrition Calibration,” preprint.
- Fan, J., Heckman, N. E. and Wand, M. P. (1995), “Local Polynomial Kernel Regression for Generalized Linear Models and Quasilikelihood Functions,” *Journal of the American Statistical Association*, 90, 141-150.
- Horvitz, D. G. and Thompson, D. J. (1952), “A Generalization of Sampling Without Replacement from a Finite Universe,” *Journal of the American Statistical Association*, 47, 663-685.
- Little, R. J. A. and Rubin, D. B. (1987). *Statistical Analysis with Missing Data*, New York: John Wiley & Sons.

- McCullagh, P. and Nelder, J. A. (1989), *Generalized Linear Models*, Second Edition, London: Chapman and Hall.
- Nelder, J. A. and Wedderburn, R. W. M. (1972), “Generalized Linear Models,” *Journal of the Royal Statistical Society, Ser. A*, 135, 370-384.
- Prentice, R. L. and Pyke, R. (1979), “Logistic Disease Incidence Models and Case-Control Studies,” *Biometrika*, 66, 403-411.
- Robins, J.M., Rotnitzky, A. and Zhao, L.P. (1994), “Estimation of Regression Coefficients When Some Regressors Are Not Always Observed,” *Journal of the American Statistical Association*, 89, 846-866.
- Rubin, D.B. (1976), “Inference and Missing Data,” *Biometrika* 63, 581-592.
- Schucany, W. R. (1995), “Adaptive Bandwidth Choice for Kernel Regression,” *Journal of the American Statistical Association*, 90, 535-540.
- Severini, T. A. and Staniswalis, J. G. (1994), “Quasilielihood Estimation in Semiparametric Models,” *Journal of the American Statistical Association*, 89, 501-511.
- Staniswalis, J. G. (1989), “The Kernel Estimate of a Regression Function in Likelihood-based Models,” *Journal of the American Statistical Association*, 84, 276-283.
- Wang, C. Y., Wang, S., Zhao, L. P. and Ou, S. T. (1997), “Weighted semiparametric estimation in regression analysis with missing covariate data,” *Journal of the American Statistical Association*, in press.
- Wedderburn, R. W. M. (1974), “Quasilielihood Functions, Generalized Linear Models, and the Gauss-Newton Method,” *Biometrika*, 61, 439-447.
- White, J. E. (1982), “A Two Stage Design for the Study of the Relationship between a Rare Exposure and a Rare Disease,” *American Journal of Epidemiology*, 115, 119-128.

## APPENDIX: TECHNICAL PROOFS

### A.1 Proof of Theorem 1

Firstly, we present a brief proof of the limit distribution of  $\hat{\mu}(\cdot, \pi)$ . The readers are referred to Fan, et al. (1995) for some related calculations. Recall that we use known  $\pi$  now. Define  $\rho(x) = ([g'\{\mu(x)\}]^2 V\{\mu(x)\})^{-1}$ , and let  $q_i(x, y) = (\partial^i)(\partial x^i)Q\{g^{-1}(x), y\}$ . Fan, et al. (1995) noted that  $q_i$  is linear in  $y$  for a fixed  $x$  and that  $q_1\{\eta(x), \mu(x)\} = 0$  and  $q_2\{\eta(x), \mu(x)\} = -\rho(x)$ .

*Conditions:*

(A1) The function  $q_2(x, y) < 0$  for  $x \in R$  and  $y$  in the range of the response variable.

(A2) The function  $f'_X, \eta^{(3)}, \text{var}(Y|X = \cdot), V^{(2)}$  and  $g^{(3)}$  are continuous.

(A3) For each  $x \in \text{supp}(f_X), \rho(x), \text{var}(Y|X = x)$ , and  $g'\{\mu(x)\}$  are nonzero.

(A4) The kernel function  $K$  is a symmetric probability density with support  $[-1, 1]$ .

(A5) For each point  $x_0$  on boundary of  $\text{supp}(f_X)$  there exists a nontrivial interval  $\mathcal{C}$  containing  $x_0$  such that  $\inf_{x \in \mathcal{C}} f_X(x) > 0$ .

(A6) The selection probability  $\pi(y) > 0$  for all  $y \in \text{supp}(f_Y)$ .

(A7)  $E[q_1\{\eta(X_1), Y_1\}(\delta_1/\pi_1)]^{2+\epsilon} < \infty$  for some  $\epsilon > 0$ .

**Proof of the Asymptotic Distribution of  $\hat{\mu}(\cdot, \pi)$ .** We study the asymptotic properties of  $\hat{\beta}^* = (nh)^{1/2}[\hat{\beta}_0 - \eta(x), h\{\hat{\beta}_1 - \eta'(x)\}]^t$ . Let  $\bar{\eta}(x, u) = \eta(x) + \eta'(x)(u - x)$ ,  $X_i^* = \{1, (X_i - x)/h\}^t$  and  $\beta^* = (nh)^{1/2}[\beta_0 - \eta(x), h\{\beta_1 - \eta'(x)\}]^t$ . Since  $\beta_0 + \beta_1(X_i - x) = \bar{\eta}(x, X_i) + (nh)^{-1/2}\beta^{*t}X_i^*$ , if  $(\hat{\beta}_0, \hat{\beta}_1)$  maximizes (3), then  $\hat{\beta}^*$  maximize

$$\sum_{i=1}^n Q[g^{-1}\{\bar{\eta}(x, X_i) + (nh)^{-1/2}\beta^{*t}X_i^*\}, Y_i] \frac{\delta_i}{\pi_i} K_h(X_i - x), \quad (8)$$

as a function of  $\beta^*$ , where  $\pi_i = \pi(Y_i)$ . We consider the normalized function

$$l_n(\beta^*, \pi) = \sum_{i=1}^n \left( Q[g^{-1}\{\bar{\eta}(x, X_i) + (nh)^{-1/2}\beta^{*t}X_i^*\}, Y_i] - Q[g^{-1}\{\bar{\eta}(x, X_i)\}, Y_i] \right) \frac{\delta_i}{\pi_i} K_h(X_i - x). \quad (9)$$

Then  $\hat{\beta}^* = \hat{\beta}^*(\pi)$  maximizes  $l_n(\cdot, \pi)$ . Let

$$W_n(\pi) = (nh)^{-1/2} \sum_{i=1}^n q_1\{\bar{\eta}(x, X_i), Y_i\} \frac{\delta_i}{\pi_i} K_h(X_i - x) X_i^*; \quad (10)$$

$$A_n(\pi) = (nh)^{-1} \sum_{i=1}^n q_2\{\bar{\eta}(x, X_i), Y_i\} \frac{\delta_i}{\pi_i} K_h(X_i - x) X_i^* X_i^{*t}. \quad (11)$$

Similar to Fan, et al. (1995), we have that

$$\begin{aligned} l_n(\beta^*, \pi) &= W_n^t(\pi) \beta^* + \frac{1}{2} \beta^{*t} A_n(\pi) \beta^* + O_p\{(nh)^{-1/2}\} \\ &= W_n^t(\pi) \beta^* - \frac{1}{2} \beta^{*t} (\Sigma_x + h\Lambda_x) \beta^* + O_p\{(nh)^{-1/2}\} + o_p(h), \end{aligned}$$

where

$$\Sigma_x = \rho(x) f_X(x) \begin{pmatrix} \gamma_0(N_x^h) & \gamma_1(N_x^h) \\ \gamma_1(N_x^h) & \gamma_2(N_x^h) \end{pmatrix}; \quad \Lambda_x = (\rho f_X)'(x) \begin{pmatrix} \gamma_1(N_x^h) & \gamma_2(N_x^h) \\ \gamma_2(N_x^h) & \gamma_3(N_x^h) \end{pmatrix}. \quad (12)$$

By the Quadratic Approximation Lemma of Fan, et al. (1995) and under the bandwidth condition that  $nh^3 \rightarrow \infty$ , we have that

$$\hat{\beta}^* = \Sigma_x^{-1} W_n(\pi) - h \Sigma_x^{-1} \Lambda_x \Sigma_x^{-1} W_n(\pi) + o_p(h), \quad (13)$$

Similar to Fan, et al. (1995), we can show that

$$\begin{aligned} E\{W_n(\pi)\} &= \frac{1}{2} (nh^5)^{1/2} \eta^{(2)}(x) \rho(x) f_X(x) \begin{pmatrix} \gamma_2(N_x^h) \\ \gamma_3(N_x^h) \end{pmatrix} + O\{(nh^7)^{1/2}\} \\ &\equiv n^{1/2} h^{5/2} B_x + O\{(nh^7)^{1/2}\}; \\ \text{var}\{W_n(\pi)\} &= \frac{f_X(x) \mathcal{L}(x)}{[V\{\mu(x)\} g'\{\mu(x)\}]^2} \begin{pmatrix} \tau_0(N_x^h) & \tau_1(N_x^h) \\ \tau_1(N_x^h) & \tau_2(N_x^h) \end{pmatrix} + o(h) \equiv \Gamma_x + o(h), \end{aligned} \quad (14)$$

where  $\mathcal{L}(x)$  is given in (4).

It can be shown by checking the Lyapounov's condition and using the Cramér-Wold device that  $\hat{\beta}^*$  is asymptotically normally distributed. From (13), we get the approximations

$$E(\hat{\beta}^*) = \Sigma_x^{-1} n^{1/2} h^{5/2} B_x + O\{(nh^7)^{1/2}\} + o(h); \quad \text{var}(\hat{\beta}^*) = \Sigma_x^{-1} \Gamma_x \Sigma_x^{-1} + o(h).$$

The proof of the first part of Theorem 1 thus follows since we are only concerned with the first component of  $\hat{\beta}^*$ , and  $\mu(x) = g^{-1}\{\eta(x)\}$ .

We now present some additional conditions for dealing with the asymptotic distribution of  $\hat{\mu}(\cdot, \hat{\pi})$ . Define  $\rho^*(y) = [\{g^{*(1)}(\pi(y))\}^2 V^*\{\pi(y)\}]^{-1}$ , and let  $q_i^*(y, z) = (\partial^i)(\partial y^i) Q^*(g^{*-1}(y), z)$ . Again, we have that

$$q_1^*\{\eta^*(y), \pi(y)\} = 0 \quad \text{and} \quad q_2^*\{\eta^*(y), \pi(y)\} = -\rho^*(y). \quad (15)$$

In addition to Condition (A1)-(A7), we need the following conditions.

*Conditions:*

- (B1) The function  $q_2^*(y, \delta) < 0$  for  $y \in R$  and  $\delta = 0, 1$ .
- (B2) The function  $f_Y', \eta^{*(3)}, \text{var}(\delta|Y = \cdot), V^{*(2)}, g^{*(3)}$  and  $\pi^{(2)}$  are continuous.
- (B3) For each  $y \in \text{supp}(f_Y), \rho^*(x), V^*(y)$  and  $g'\{\pi(y)\}$  are nonzero.
- (B4) For each point  $y_0$  on boundary of  $\text{supp} f_Y$  there exists a nontrivial interval  $\mathcal{C}$  containing  $y_0$  such that  $\inf_{y \in \mathcal{C}} f_Y(y) > 0$ .
- (B5)  $\inf\{\pi(y) : y \in \text{supp}(f_Y)\} > 0$ .
- (B6) The conditional density of  $X$  given  $Y$  is bounded a.e.

Before proving the main part of Theorem 1, we present some lemmas which will be used in the proof. Recall that  $\hat{\pi}$  was defined in (5).

**Lemma 1.** Under the same conditions as those of Theorem 1,  $G_n = o_p(h)$ , where

$$G_n = (nh)^{-1} \sum_{i=1}^n q_2\{\bar{\eta}(x, X_i), Y_i\} K_h(X_i - x) X_i^* X_i^{*t} \frac{\delta_i}{\pi_i^2} (\hat{\pi}_i - \pi_i); \quad X_i^* = \{1, (X_i - x)/h\}^t.$$

**Lemma 2.** Under the same conditions as those in Theorem 1,  $C_n = o_p(h^{1/2})$ , where

$$C_n = (nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\delta_i - \pi_i}{\pi_i} \frac{\hat{\pi}_i - \pi_i}{\pi_i} q_1\{\bar{\eta}(x, X_i), Y_i\} K_h(X_i - x) X_i^* \right].$$

**Lemma 3.** Under the same conditions as those of Theorem 1, let

$$D_n = (nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\hat{\pi}_i - \pi_i}{\pi_i} q_1 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* \right].$$

Then there exists an  $S_3(x)$  and  $D_n^*$  with  $E(D_n^*) = o(n^{1/2}h^{5/2})$  and  $\text{var}(D_n^*) = o(h^2)$ , such that

$$D_n - n^{1/2}h^{5/2}(c^*)^2 f_X(x) S_3(x) = (nh)^{-1/2} \sum_{i=1}^n \{ (\delta_i - \pi_i) / \pi_i \} \mathcal{M}_h(Y_i) + D_n^*,$$

where

$$\mathcal{M}_h(Y_i) = E [q_1 \{ \bar{\eta}(x, X_i), Y_i \} X_i^* K_h(X_i - x) | Y_i]. \quad (16)$$

The proofs of Lemmas 1-3 will be postponed until after the proof of the limit distribution of  $\hat{\mu}(\cdot, \hat{\pi})$ .

**Proof of the Asymptotic Distribution of  $\hat{\mu}(\cdot, \hat{\pi})$ .** Recall that  $\hat{\beta} = \hat{\beta}(\hat{\pi})$  maximizes (6), and we defined  $l_n(\beta^*, \pi)$  in (9). The main step here is to derive the asymptotic expression of  $l_n(\beta^*, \hat{\pi})$ . Note that  $A_n(\pi)$  was defined in (11). We have

$$\begin{aligned} A_n(\hat{\pi}) - A_n(\pi) &= \left[ (nh)^{-1} \sum_{i=1}^n q_2 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* X_i^{*t} \frac{\delta_i}{\pi_i^2} (\hat{\pi}_i - \pi_i) \right] \{1 + o_p(1)\} \\ &\equiv G_n \{1 + o_p(1)\}. \end{aligned}$$

By Lemma 1, we have  $A_n(\hat{\pi}) = A_n(\pi) + o_p(h)$ . Using similar calculations as in the proof of the distribution of  $\hat{\mu}(\cdot, \pi)$ ,

$$l_n(\beta^*, \hat{\pi}) = W_n^t(\hat{\pi}) \beta^* - \frac{1}{2} \beta^{*t} (\Sigma_x + h \Lambda_x) \beta^* + O_p\{(nh)^{-1/2}\} + o_p(h).$$

For simplicity in this proof we continue to use  $\hat{\beta}^* = \hat{\beta}^*(\hat{\pi})$  as the maximizer of (8) with estimated  $\hat{\pi}$ . By the Quadratic Approximation Lemma of Fan, et al. (1995), we have that

$$\hat{\beta}^* = \Sigma_x^{-1} W_n(\hat{\pi}) - h \Sigma_x^{-1} \Lambda_x \Sigma_x^{-1} W_n(\pi) + o_p(h). \quad (17)$$

We now find the limit distribution of  $W_n(\hat{\pi})$ , where  $W_n(\cdot)$  was defined in (10). By a linearization technique as in Wang, et al. (1997), we have

$$\begin{aligned} W_n(\hat{\pi}) &= (nh)^{-1/2} \sum_{i=1}^n \left[ q_1 \{ \bar{\eta}(x, X_i), Y_i \} \frac{\delta_i}{\pi_i} K_h(X_i - x) X_i^* \left( 1 - \frac{\hat{\pi}_i - \pi_i}{\pi_i} \right) \right] \\ &\quad + \left( (nh)^{-1/2} \sum_{i=1}^n \left[ q_1 \{ \bar{\eta}(x, X_i), Y_i \} \frac{\delta_i}{\pi_i} K_h(X_i - x) X_i^* \frac{(\hat{\pi}_i - \pi_i)^2}{\pi_i} \right] \right) \{1 + o_p(1)\}. \end{aligned}$$

Denote the second term of the above equation by  $R_n$ . Then it can be shown to have mean  $O(n^{1/2}h^{9/2})$  and variance  $o(h)$ . Therefore, we have that

$$\begin{aligned} W_n(\hat{\pi}) &= (nh)^{-1/2} \sum_{i=1}^n \left( \frac{\delta_i}{\pi_i} - \frac{\hat{\pi}_i - \pi_i}{\pi_i} \right) q_1 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* \\ &\quad - (nh)^{-1/2} \sum_{i=1}^n \frac{\delta_i - \pi_i}{\pi_i} \frac{\hat{\pi}_i - \pi_i}{\pi_i} q_1 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* + R_n \\ &\equiv W_n(\pi) - D_n - C_n + R_n. \end{aligned} \quad (18)$$

By Lemmas 2 and 3, we have that

$$W_n(\hat{\pi}) = W_n(\pi) - n^{1/2} h^{5/2} (c^*)^2 f_X(x) S_3(x) - (nh)^{-1/2} \sum_{i=1}^n \frac{\delta_i - \pi_i}{\pi_i} \mathcal{M}_h(Y_i) + R_n^* \quad (19)$$

for some  $R_n^*$  that has mean  $o(n^{1/2}h^{5/2})$  and variance  $o(1)$ . Let  $f_{X|Y}$  denote the conditional density of  $X$  given  $Y$ . By direct calculations,  $\mathcal{M}_h(Y_j) = h \left[ q_1 \{ \eta(x), Y_j \} f_{X|Y}(x) \{ \gamma_0(N_x^h), \gamma_1(N_x^h) \}^t \right] \{ 1 + o_p(1) \} \equiv h \mathcal{D}(Y_j) \{ 1 + o_p(1) \}$ . By (14), (17) and (19), the asymptotic distribution of  $\hat{\mu}(\cdot, \hat{\pi})$  follows because

$$\begin{aligned} E \left\{ \frac{\delta_i - \pi_i}{\pi_i} \mathcal{M}_h(Y_i) \right\} &= E \left[ E \left\{ \frac{\delta_i - \pi_i}{\pi_i} \mathcal{M}_h(Y_i) | Y_i \right\} \right] = 0; \\ \text{var} \left\{ (nh)^{-1/2} \sum_{i=1}^n \frac{\delta_i - \pi_i}{\pi_i} \mathcal{M}_h(Y_i) \right\} &= h^{-1} E \left\{ \mathcal{M}_h(Y_1) \mathcal{M}_h^t(Y_1) \text{var} \left( \frac{\delta_1 - \pi_1}{\pi_1} | Y_1 \right) \right\} \\ &= h^{-1} E \left\{ \frac{1 - \pi_1}{\pi_1} \mathcal{M}_h(Y_1) \mathcal{M}_h^t(Y_1) \right\} = h E \left\{ \frac{1 - \pi_1}{\pi_1} \mathcal{D}(Y_1) \mathcal{D}^t(Y_1) \right\} \{ 1 + o(1) \} = O(h). \end{aligned}$$

The last equation holds by Conditions (B5) and (B6).

**Proof of Lemma 1.** Firstly, we note that  $E\{(\hat{\pi}_1 - \pi_1) | Y_1\} = \lambda^2 c_1(Y_1) \{ 1 + o_p(1) \}$  and  $\text{var}\{(\hat{\pi}_1 - \pi_1) | Y_1\} = (n\lambda)^{-1} c_2(Y_1) \{ 1 + o_p(1) \}$ , for some functions  $c_1$  and  $c_2$ . Let  $\hat{\pi}_{i(j)}$  denote  $\hat{\pi}_i$  without using subject  $j$ . Then

$$\begin{aligned} E(G_n) &= E[q_2 \{ \bar{\eta}(x, X_1), Y_1 \} \frac{1}{h} K_h(X_1 - x) X_1^* X_1^{*t} \frac{\delta_1}{\pi_1^2} (\hat{\pi}_1 - \pi_1)] \\ &= E \left( E[q_2 \{ \bar{\eta}(x, X_1), Y_1 \} \frac{1}{h} K_h(X_1 - x) X_1^* X_1^{*t} \frac{\delta_1}{\pi_1^2} (\hat{\pi}_{1(1)} - \pi_1) | X_1, Y_1] \right) + O\left(\frac{1}{n\lambda}\right) \\ &= E[q_2 \{ \bar{\eta}(x, X_1), Y_1 \} \frac{1}{h} K_h(X_1 - x) X_1^* X_1^{*t} \frac{1}{\pi_1} E\{(\hat{\pi}_1 - \pi_1) | Y_1\}] + O\left(\frac{1}{n\lambda}\right) \\ &= \lambda^2 E[q_2 \{ \bar{\eta}(x, X_1), Y_1 \} \frac{1}{h} K_h(X_1 - x) X_1^* X_1^{*t} \frac{1}{\pi_1} c_1(Y_1) \{ 1 + o_p(1) \}] + O\left(\frac{1}{n\lambda}\right) = o(h). \end{aligned}$$

The last equation holds since  $\lambda = c^* h$  for some  $c^* > 0$ . Note that if we let  $S_i = q_2 \{ \bar{\eta}(x, X_i), Y_i \} h^{-1} K_h(X_i - x) X_i^* X_i^{*t} (\delta_i / \pi_i^2) (\hat{\pi}_i - \pi_i)$ , then following calculations similar to those above we obtain that

$\text{cov}(S_i, S_j) = o(h^2)$  when  $i \neq j$ . Therefore, the variance of the left-upper element of  $G_n$  is

$$\begin{aligned}
\text{var}[\{G_n\}_{11}] &= n^{-1} \text{var}[q_2\{\bar{\eta}(x, X_1), Y_1\} \frac{1}{h} K_h(X_1 - x) \frac{\delta_1}{\pi_1^2} (\hat{\pi}_1 - \pi_1)] + o(h^2) \\
&= n^{-1} E[q_2^2\{\bar{\eta}(x, X_1), Y_1\} \frac{1}{h^2} K_h^2(X_1 - x) \frac{1}{\pi_1^4} \text{var}\{(\hat{\pi}_{1(1)} - \pi_1)|Y_1\}] \\
&\quad + n^{-1} \text{var}[q_2\{\bar{\eta}(x, X_1), Y_1\} \frac{1}{h} K_h(X_1 - x) \frac{\delta_1}{\pi_1^2} E\{(\hat{\pi}_{1(1)} - \pi_1)|Y_1\}] + o(h^2) \\
&= (nh)^{-1} E[q_2^2\{\bar{\eta}(x, X_1), Y_1\} \frac{1}{h} K_h^2(X_1 - x) \frac{1}{\pi_1^3} (n\lambda)^{-1} c_2(Y_1) \{1 + o_p(1)\}] \\
&\quad + (nh)^{-1} \text{var}[q_2\{\bar{\eta}(x, X_1), Y_1\} \frac{1}{\sqrt{h}} K_h(X_1 - x) \frac{\delta_1}{\pi_1^2} \lambda^2 c_1(Y_1) \{1 + o_p(1)\}] + o(h^2) \\
&= O\{(nh)^{-1} (n\lambda)^{-1} + (nh)^{-1} \lambda^4\} + o(h^2) = o(h^2).
\end{aligned}$$

The last equation holds since  $\lambda = c^*h$  and  $nh^3 \rightarrow \infty$ . Similar calculations lead to  $\text{var}[\{G_n\}_{12}] = o(h^2)$  and  $\text{var}[\{G_n\}_{22}] = o(h^2)$ , completing the proof of Lemma 1.

**Proof of Lemma 2.** We assume that  $Y$  is a continuous random variable; a similar approach can be easily applied to discrete  $Y$ . Using calculations similar to Fan, et al. (1995) and under Conditions (B1)-(B4), we can show that for each  $y$ , as  $n\lambda \rightarrow \infty$ ,

$$\hat{\pi}(y) - \pi(y) = (n\lambda)^{-1/2} [g^{*(1)}\{\pi(y)\}]^{-1} \{\Sigma_y^{*-1} W_n^*\}_1 + (n\lambda)^{-1/2} \mathcal{R}(y, \tilde{\delta}, \tilde{Y}), \quad (20)$$

where  $\{\cdot\}_1$  denotes the first component of a vector,  $\tilde{\delta} = (\delta_1, \dots, \delta_n)$ ,  $\tilde{Y} = (Y_1, \dots, Y_n)$ ,

$$W_n^* = (n\lambda)^{-1/2} \sum_{i=1}^n q_1^* \{\bar{\eta}^*(y, Y_i), \delta_i\} K_\lambda(Y_i - y) Y_i^*; \quad \Sigma_y^* = \rho^*(y) f_Y(y) \begin{pmatrix} \gamma_0(N_y^\lambda) & \gamma_1(N_y^\lambda) \\ \gamma_1(N_y^\lambda) & \gamma_2(N_y^\lambda) \end{pmatrix};$$

$\bar{\eta}^*(y, u) = \eta^*(y) + \eta^{*(1)}(u)(u - y)$ ,  $Y_i^* = (1, (Y_i - y)/\lambda)^t$ ,  $N_y^\lambda = \{z : y - \lambda z \in \text{supp}(f_Y) \cap [-1, 1]\}$ ,  $\mathcal{R}(y, \tilde{\delta}, \tilde{Y}) = -\lambda \left[ g^{*(1)}\{\pi(y)\} \right]^{-1} \{\Sigma_y^{*(-1)} \Lambda_y^* \Sigma_y^{*(-1)} W_n^*\}_1 \{1 + o_p(1)\}$ ,  $\Lambda_y^*$  is the same as  $\Lambda_x$  defined in (12) except replacing  $\rho f_X$  by  $\rho^* f_Y$ ,  $f_Y(\cdot)$  being the density of  $Y$ . From (20), we have that there is a function  $g_\pi^*(y)$  and a function  $\mathcal{X}(\delta, y)$  such that

$$\hat{\pi}(y) - \pi(y) = \lambda^2 g_\pi^*(y) + (n\lambda)^{-1} \sum_{i=1}^n K_\lambda(Y_i - y) \mathcal{X}(\delta_i, Y_i) + o_p(\lambda^2) + o_p(n^{-1/2}), \quad (21)$$

where  $E\{\mathcal{X}(\delta_i, Y_i)|Y_i\} = 0$  and the  $o_p(\lambda^2)$  term does not depend on  $\delta$ 's. Note that as in Fan, et al. (1995), it can be shown that the term associated with the bias  $g_\pi^*(y) = 0$  if  $\pi'(y) = 0$ . Therefore,

$$(nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\delta_i - \pi_i}{\pi_i} \frac{\lambda^2 g_\pi^*(Y_i)}{\pi_i} q_1\{\bar{\eta}(x, X_i), Y_i\} K_h(X_i - x) X_i^* \right] = O_p(\lambda^2)$$

by calculating the mean and the variance. Also,

$$(nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\delta_i - \pi_i}{\pi_i^2} \left\{ (n\lambda)^{-1} \sum_{j=1}^n K_\lambda(Y_j - Y_i) \mathcal{X}(\delta_j, Y_j) \right\} q_1\{\bar{\eta}(x, X_i), Y_i\} K_h(X_i - x) X_i^* \right]$$



can be shown to be  $o_p(h)$  because each summand has a factor  $(\delta_i - \pi_i)\mathcal{X}(\delta_j, Y_j)$  which has mean 0 if  $i \neq j$ . Further,

$$(nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\delta_i - \pi_i}{\pi_i^2} q_1 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* \{ o_p(\lambda^2) + o_p(n^{-1/2}) \} \right] = o_p(h^{1/2}),$$

because the  $o_p(\lambda^2)$  term above does not depend on the  $\delta$ 's. Therefore  $C_n = o_p(h^{1/2})$ .

**Proof of Lemma 3.** Applying (21) to  $D_n$  and write  $D_n = D_{1n} + D_{2n} + D_{3n}$ . First,

$$D_{1n} = (nh)^{-1/2} \sum_{i=1}^n \left[ \frac{\lambda^2 g_\pi^*(Y_i)}{\pi_i} q_1 \{ \bar{\eta}(x, X_i), Y_i \} K_h(X_i - x) X_i^* \right].$$

Then note that

$$\begin{aligned} E(D_{1n}) &= n^{1/2} h^{-1/2} \lambda^2 \int \frac{g_\pi^*(y_1)}{\pi_1} q_1 \{ \bar{\eta}(x, x_1), y_1 \} K_h(x_1 - x) x_1^* f_{Y,X}(y_1, x_1) dy_1 dx_1 \\ &= n^{1/2} h^{-1/2} \lambda^2 \int \frac{g_\pi^*(y_1)}{\pi_1} \frac{y_1 - \mu(x) - [g^{(1)}\{\mu(x)\}]^{-1} \eta^{(1)}(x)(x_1 - x)}{g^{(1)}\{\mu(x)\} V\{\mu(x)\}} K_h(x_1 - x) x_1^* \\ &\quad \times f_{Y,X}(y_1, x_1) dy_1 dx_1 \end{aligned}$$

Let

$$\begin{aligned} S_1(x) &= E \left( g_\pi^*(Y_1) Y_1 [\pi(Y_1) g^{(1)}\{\mu(x)\} V\{\mu(x)\}]^{-1} \right); \\ S_2(x) &= E \left( g_\pi^*(Y_1) [\pi(Y_1) g^{(1)}\{\mu(x)\} V\{\mu(x)\}]^{-1} \right), \end{aligned}$$

and

$$S_3(x) = S_1(x) - \mu(x) S_2(x). \quad (22)$$

Note that  $S_3(x) = 0$  if either  $Y$  is a lattice random variable or  $\pi'(Y) = 0$  *a.e.*, because under these circumstances  $g_\pi^*(Y) = 0$  *a.e.* Then it is easily seen that  $E\{(D_{1n})_1\} = n^{1/2} h^{5/2} (c^*)^2 f_X(x) S_3(x) + o_p(n^{1/2} h^{5/2})$  and it can also be shown that  $\text{var}\{(D_{1n})_1\} = O(\lambda^4)$ .

Now consider

$$D_{2n} = (nh)^{-1/2} \sum_{j=1}^n \mathcal{X}(\delta_j, Y_j) \left( (n\lambda)^{-1} \sum_{i=1}^n \left[ \frac{q_1 \{ \bar{\eta}(x, X_i), Y_i \} X_i^*}{\pi_i} K_h(X_i - x) \right] K_\lambda(Y_j - Y_i) \right).$$

To estimate  $D_{2n}$ , we note that by some further calculations the  $\mathcal{X}(\delta_i, Y_i)$  in (21) can be written as  $\mathcal{X}(\delta_i, Y_i) = \{f_Y(y)\psi(y)\}^{-1} \phi(Y_i, y) \{\delta_i - \pi(Y_i)\} + o_p(1)$ , where  $\phi(Y_i, y) = \gamma_2(N_y^\lambda) - \gamma_1(N_y^\lambda) \{(Y_i - y)/\lambda\}$  and  $\phi(y) = \gamma_0(N_y^\lambda) \gamma_2(N_y^\lambda) - \gamma_1^2(N_y^\lambda)$ . Therefore, by some standard calculations, we have that

$$D_{2n} = (nh)^{-1/2} \sum_{j=1}^n (\delta_j - \pi_j) \mathcal{M}_h(Y_j) / \pi_j + o_p(n^{-1/2}).$$

Finally,

$$D_{3n} = (nh)^{-1/2} \sum_{i=1}^n \left[ \frac{q_1 \{\bar{\eta}(x, X_i), Y_i\} K_h(X_i - x) X_i^*}{\pi_i} \{o_p(h^2) + o_p(n^{-1/2})\} \right].$$

By some standard calculations, we have that  $E(D_{3n}) = o(n^{1/2}h^{5/2})$ ,  $\text{var}(D_{3n}) = o(h^2)$ . Combining the calculations of  $D_{1n}, D_{2n}$  and  $D_{3n}$ , we have that there is a  $D_n^*$  such that  $E(D_n^*) = o(n^{1/2}h^{5/2})$ ,  $\text{var}(D_n^*) = o(h^2)$  and

$$D_n - n^{1/2}h^{5/2}(c^*)^2 f_X(x) S_3(x) = (nh)^{-1/2} \sum_{j=1}^n \frac{\delta_j - \pi_j}{\pi_j} \mathcal{M}_h(Y_j) + D_n^*.$$

## A.2 Proof of Theorem 2.

By (18),  $W_n(\hat{\pi}) = W_n(\pi) - D_n - C_n + R_n$ . As in the proof of Lemma 3,  $D_n = (nh)^{-1/2} \sum_{i=1}^n \{(\delta_i - \pi_i)/\pi_i\} \mathcal{M}_h(Y_i) + n^{1/2}h^{5/2}(c^*)^2 f_X(x) S_3(x) + D_n^*$ . By some standard calculations, we have

$$\begin{aligned} & \text{cov}\{W_n(\pi), D_n\} \\ &= (nh)^{-1} \sum_{i=1}^n \text{cov} \left[ q_1 \{\bar{\eta}(x, X_i), Y_i\} \frac{\delta_i}{\pi_i} K_h(X_i - x) X_i^*, \frac{\delta_i - \pi_i}{\pi_i} \mathcal{M}_h(Y_i) \right] + o(h) \\ &= h^{-1} E \left[ \frac{1 - \pi_1}{\pi_1} q_1 \{\bar{\eta}(x, X_1), Y_1\} X_1^* K_h(X_1 - x) \mathcal{M}_h^t(Y_1) \right] + o(h) \\ &= h^{-1} E \left\{ \frac{1 - \pi_1}{\pi_1} \mathcal{M}_h(Y_1) \mathcal{M}_h^t(Y_1) \right\} + o(h) \\ &= h E \left\{ \frac{1 - \pi_1}{\pi_1} \mathcal{D}(Y_1) \mathcal{D}^t(Y_1) \right\} + o(h). \end{aligned}$$

The  $\mathcal{D}(Y)$  in the above calculations was defined in the proof of Theorem 1. The covariances due to  $\text{cov}\{W_n(\pi), C_n\}$  and  $\text{cov}\{W_n(\pi), R_n\}$  can be shown to be smaller than the rate of  $\text{cov}\{W_n(\pi), D_n\}$  since  $\text{cov}(C_n) = o(h)$ ,  $\text{cov}(R_n) = o(h)$ . We may in fact apply (21) to get more precise rates, e.g.,  $\text{cov}\{W_n(\pi), C_n\} = o(h^2)$ . Therefore, letting  $\hat{B}^* = \Sigma_x^{-1} \{W_n(\pi) - D_n\} - h \Sigma_x^{-1} \Lambda_x \Sigma_x^{-1} W_n(\pi) = \hat{\beta}^*(\pi) - \Sigma_x^{-1} D_n + o_p(h) = \hat{\beta}^*(\hat{\pi}) + o_p(h)$  by (13) and (17), Theorem 2 follows since  $\text{var}\{\hat{B}^*\} = \text{var}\{\hat{\beta}^*(\pi)\} - h \Sigma_x^{-1} E[\{(1 - \pi_1)/\pi_1\} \mathcal{D}(Y_1) \mathcal{D}^t(Y_1)] \Sigma_x^{-1} + o(h)$ .

## A.3 Sketch of the Proof of Generalizations in Section 7

Here we sketch the arguments of Section 7. We renormalize so that  $\gamma_0 = \gamma_2 = 1, \gamma_1 = \gamma_3 = 0$ . Carroll, Ruppert and Welsh (1996) showed that there is a function  $g_\pi(z)$  which does not depend on the density  $f_Z(\cdot)$  of  $Z$ , and a function  $\Omega(\cdot)$  such that

$$\hat{\pi}(z) - \pi(z) = (h^2/2) g_\pi(z) + (n\lambda)^{-1} \sum_{i=1}^n K_\lambda(Z_i - z) \Omega(\tilde{Y}_i, X_i, Z_i) + o_p(h^2) + o_p(n^{-1/2}); \quad (23)$$

$$E\{\Omega(\tilde{Y}, X, Z) | Z\} = 0. \quad (24)$$

By a Taylor series expansion of (7), it is easily seen that

$$-f_X(x)\Sigma_1(x)\{\hat{\eta}(x) - \eta(x)\} \approx B_{n1} - B_{n2} + B_{n3},$$

where if  $\Psi_\pi = \frac{\partial}{\partial v}\Psi(\tilde{Y}, v, \eta); \Psi_\eta = \frac{\partial}{\partial v}\Psi(\tilde{Y}, \pi, v)$ . Then

$$\begin{aligned} B_{n1} &= (nh)^{-1} \sum_{i=1}^n K_h(X_i - x) \Psi\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\}; \\ B_{n2} &= (nh)^{-1} \sum_{i=1}^n K_h(X_i - x) \left[ \Psi\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\} - \Psi\{\tilde{Y}_i, \pi(Z_i), \eta(x) + \eta'(x)(X_i - x)\} \right]; \\ B_{n3} &= (nh)^{-1} \sum_{i=1}^n K_h(X_i - x) \Psi_\pi\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\} \{\hat{\pi}(Z_i) - \pi(Z_i)\}; \\ \Sigma_1(x) &= E[\Psi_\eta\{\tilde{Y}, \pi(Z), \eta(x)\} | X = x]. \end{aligned}$$

It is easily seen that  $B_{n2} = (h^2/2)f_X(x)\Sigma_1(x)\eta^{(2)}(x)\{1 + o_p(1)\}$ . Writing  $\Psi_{\pi_i} = \Psi_\pi\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\}$ , and writing  $\Omega_i$  similarly, we note that  $B_{n3} \approx B_{n31} + B_{n32}$ , where using (23),

$$B_{n31} = (h^2/2)(nh)^{-1} \sum_{i=1}^n K_h(X_i - x) \Psi_{\pi_i} g_\pi(Z_i) = (h^2/2)f_X(x)E\{\Psi_\pi(\cdot)g_\pi(\cdot) | X = x\}\{1 + o_p(1)\}$$

and

$$B_{n32} = n^{-2}(h\lambda)^{-1} \sum_{i=1}^n \sum_{j=1}^n K_h(X_i - x) \Psi_{\pi_i} K_\lambda(Z_j - Z_i) \Omega_j.$$

Thus, we have shown that

$$\hat{\eta}(x) - \eta(x) \approx (h^2/2) \left[ \eta^{(2)}(x) - E\{\Psi_\pi g_\pi(\cdot) | X = x\} / \Sigma_1(x) \right] - \{f_X(x)\Sigma_1(x)\}^{-1} (B_{n1} + B_{n32}). \quad (25)$$

The first term in (25) is the bias, which is independent of the design density but affected by estimation of  $\pi(\cdot)$ , as claimed. To complete the argument, we merely need to show that  $B_{n32} = o_p\{(nh)^{-1/2}\}$ . Recalling that  $K(\cdot)$  is symmetric, rewrite

$$B_{n32} = n^{-1} \sum_{i=1}^n \Omega_i \{ n^{-1} \sum_{j=1}^n h^{-1} K_h(X_j - x) \lambda^{-1} K_\lambda(Z_j - Z_i) \Psi_{\pi_j} \}. \quad (26)$$

Using Chebychev's inequality, detailed algebra gives the designed result. In the interests of space we forego the calculations, but note that the term in brackets in (26) is a *bivariate* kernel regression of  $\Psi_\pi\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\}$  on  $(X, Z)$  evaluated at  $X = x, Z = Z_i$ , and hence converges to  $r(x, Z_i)$ , where  $r(x, Z_i) = E[\Psi_\pi\{\tilde{Y}, \pi(Z), \eta(X)\} | X = x, Z = Z_i]$ . Therefore,  $B_{n32} \approx n^{-1} \sum_{i=1}^n \Omega\{\tilde{Y}_i, \pi(Z_i), \eta(X_i)\} r(x, Z_i)$  which is  $O_p(n^{-1/2})$  from (24).

#### A.4 Bandwidth Selection

### Bandwidth for the Selection Probability Estimation

Fan, et al. (1995) suggested a bandwidth selector based on “plugging-in” estimates of unknown quantities. For the rest of the paper, the notation  $\phi^{(k)}(\cdot)$  denotes the  $k$ th derivative of a function  $\phi(\cdot)$ . Because we consider the local linear smoother for  $\pi$ , an approximate asymptotic mean integrated square error for  $\hat{\eta}^*(\lambda)$  is

$$\text{AMISE}\{\hat{\eta}^*(\lambda)\} = \frac{\lambda^4 \gamma_2^2}{4} \int \{\eta^{*(2)}(y)\}^2 f_Y(y) dy + (n\lambda)^{-1} \tau_0 \int V^* \{\pi(y)\} [g^{*'}\{\pi(y)\}]^2 dy,$$

where  $\gamma_2 = \gamma_2([-1, 1])$  and  $\tau_0 = \tau_0([-1, 1])$  are given in Section 3 and  $f_Y(y)$  denotes the density of  $Y$ . With respect to this criteria, the optimal bandwidth for the estimate of  $\pi$  is then

$$\lambda_{AMISE} = \left[ \frac{\tau_0 \int V \{\pi(y)\} [g^{*'}\{\pi(y)\}]^2 dy}{n \gamma_2^2 \int \{\eta^{*(2)}(y)\}^2 f_Y(y) dy} \right]^{1/5}.$$

Note that  $\pi(y)$  and  $f_Y(y)$  are unknown. An “ad-hoc” plug-in bandwidth selection is to estimate  $\eta^*(y)$  by a 3rd (or higher) degree polynomial parametric fit to the selection probabilities and to estimate  $f_Y(y)$  by a usual kernel estimate. We also note that this criteria is an approximation which does not consider the  $\gamma_0$  and  $\tau_0$  as a function of  $\lambda$  on the boundary points. In practice this selector is reasonable well for a wide range of functions.

### Bandwidth for the Primary Estimation

Now we study the bandwidth selection for our primary estimation. An approximate asymptotic mean integrated square error for  $\hat{\eta}$  is

$$\text{AMISE}\{\hat{\eta}(h)\} = \frac{h^4 \gamma_2^2}{4} \int \{\eta^{(2)}(x)\}^2 f_X(x) dy + (nh)^{-1} \tau_0 \int \mathcal{L}(x) [g'\{\mu(x)\}]^2 dx,$$

With respect to this criterion, the optimal bandwidth for the estimate of  $\mu$  is then

$$h_{AMISE} = \left[ \frac{\tau_0 \int \mathcal{L}(x) [g'\{\mu(x)\}]^2 dx}{n \gamma_2^2 \int \{\eta^{(2)}(x)\}^2 f_X(x) dx} \right]^{1/5}.$$

Similar to the argument of the selection of  $\lambda$ , we may estimate  $\eta(x)$  by a third (or higher) degree polynomial. In addition, we may estimate  $\mathcal{L}(x) = E[\{Y_1 - \mu(X_1)\}^2 / \pi(Y_1) | X_1 = x]$  and  $f_X(x)$  by nonparametric estimation based on validation data with inverse selection weights. This gives a global bandwidth selection. Alternatively, Schucany (1995) proposed an adaptive local bandwidth estimator for the Nadaraya-Watson estimator, and found that it has improvements over a global bandwidth estimator. It maybe a worthwhile future project to study the local bandwidth selector in the problem of generalized linear missing data models.

C. Y. Wang  
Division of Public Health Sciences  
Fred Hutchinson Cancer Research Center  
1124 Columbia Street, MP 1002  
Seattle, WA 98104

Suojin Wang  
Department of Statistics  
Texas A&M University  
College Station, TX 77843-3143

R. J. Carroll  
Department of Statistics  
Texas A&M University  
College Station, TX 77843-3143

Roberto G. Gutierrez  
Department of Statistical Science  
Southern Methodist University  
Dallas, TX 77275-0332

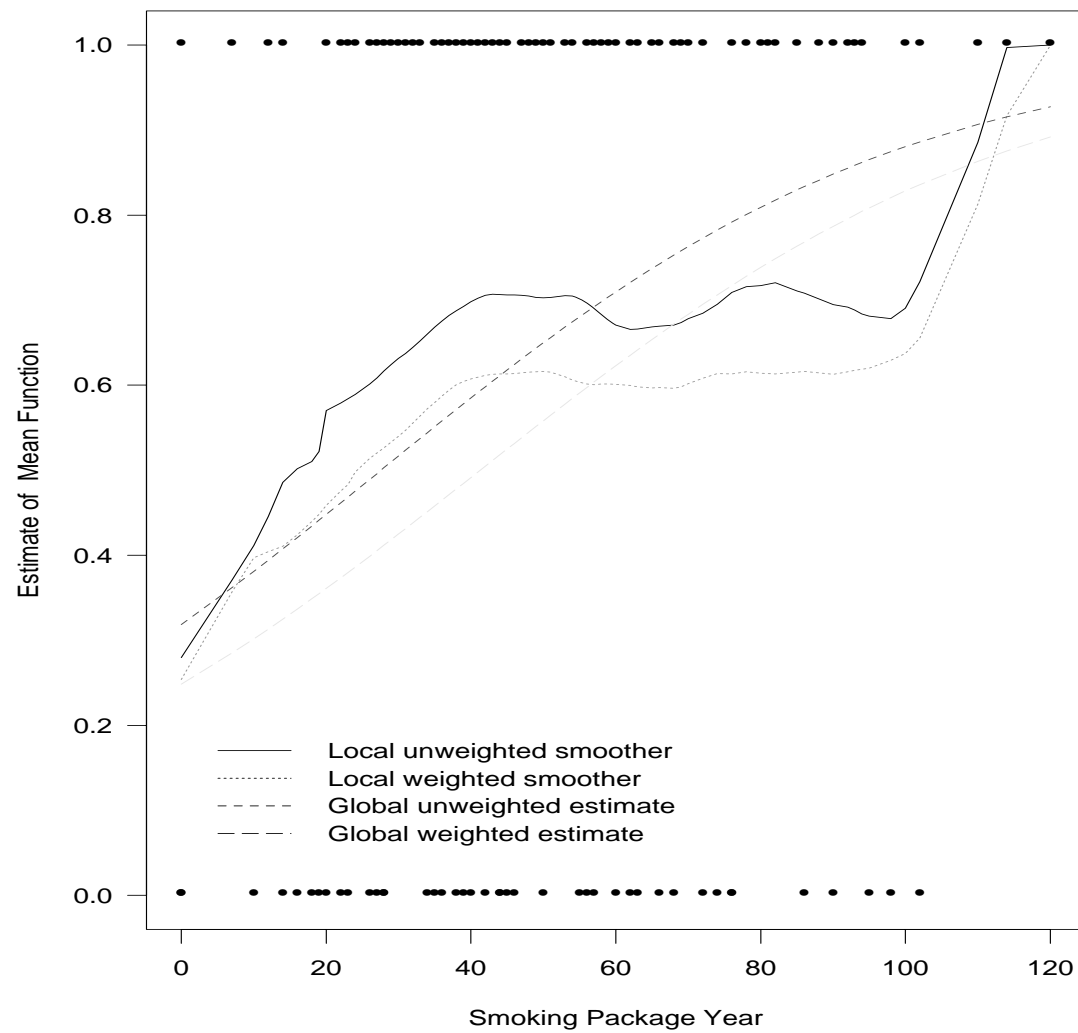


Figure 1: *Bladder Cancer Case-control Data Analysis.*

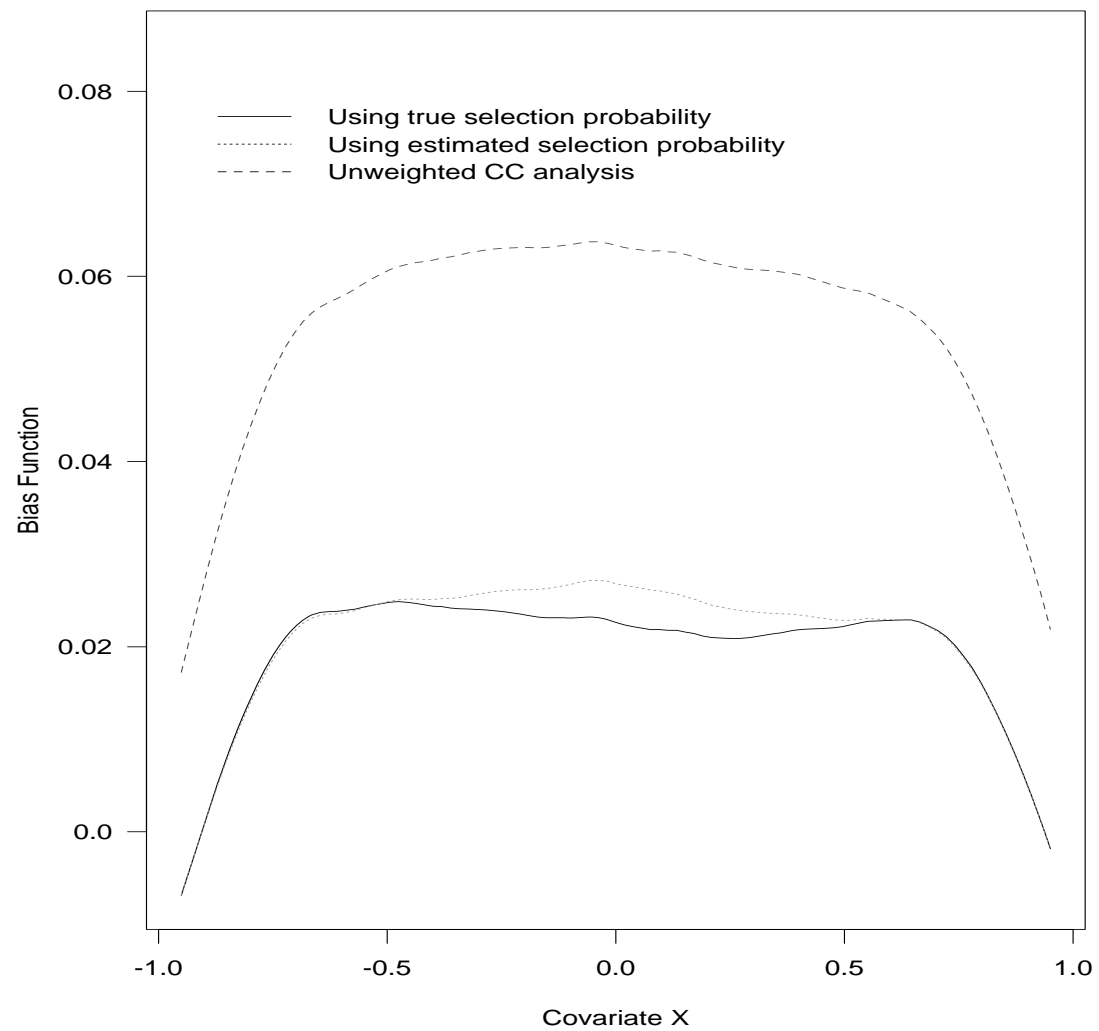


Figure 2: *Simulation Study for Biases from Estimating  $\mu$  for Continuous Response.*

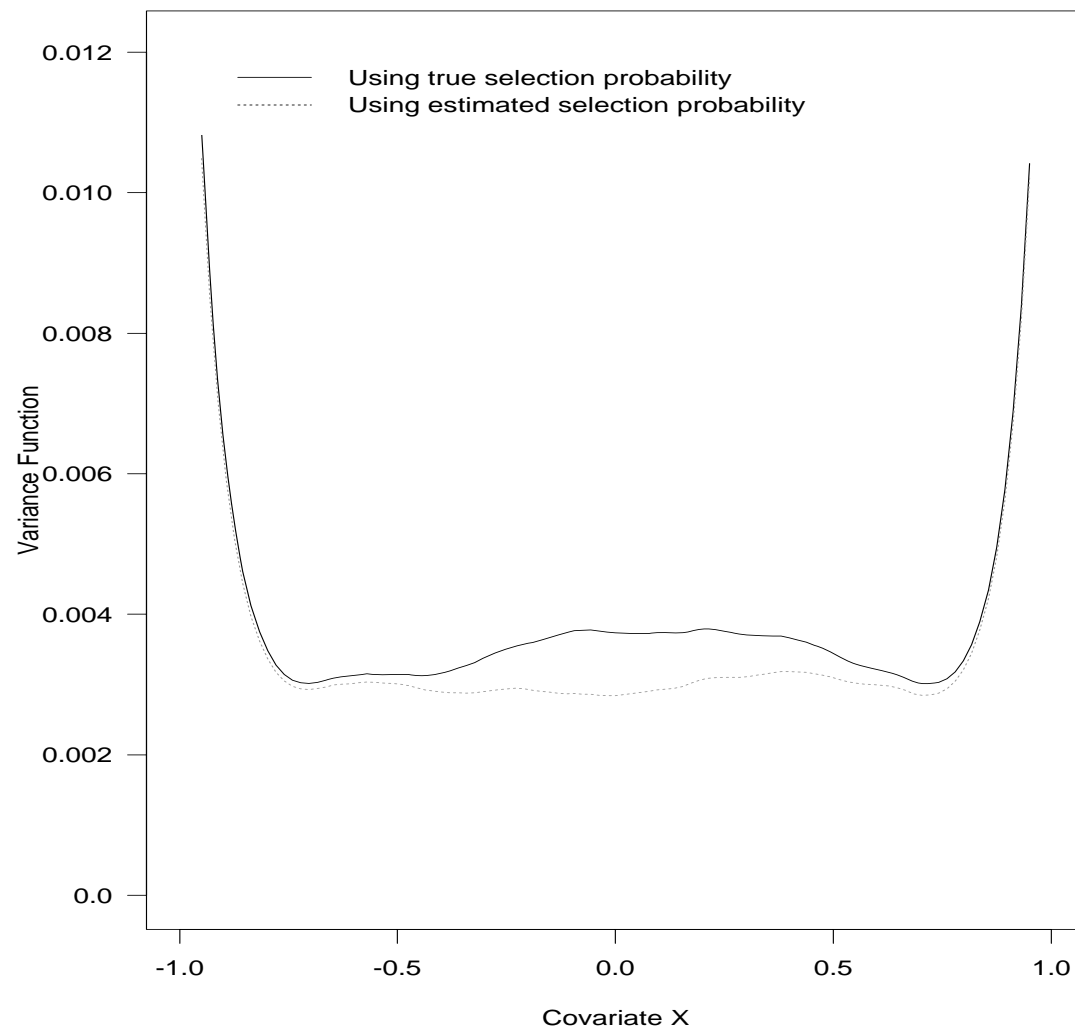


Figure 3: *Simulation Study for Variances from Estimating  $\mu$  for Continuous Response.*